

A7

6/5/1 (Item 1 from file: 347)  
DIALOG(R) File 347:JAPIO  
(c) 2000 JPO & JAPIO. All rts. reserv.

04384503 \*\*Image available\*\*  
DOCUMENT RETRIEVING DEVICE

PUB. NO.: 06-028403 JP 6028403 A]  
PUBLISHED: February 04, 1994 (19940204)  
INVENTOR(s): SASAKI MIKIRO  
ISHIKAWA HIROMICHI  
APPLICANT(s): MITSUBISHI ELECTRIC CORP [000601] (A Japanese Company or  
Corporation), JP (Japan)  
APPL. NO.: 04-182438 [JP 92182438]  
FILED: July 09, 1992 (19920709)  
INTL CLASS: [5] G06F-015/40  
JAPIO CLASS: 45.4 (INFORMATION PROCESSING -- Computer Applications)  
JAPIO KEYWORD: R139 (INFORMATION PROCESSING -- Word Processors)  
JOURNAL: Section: P, Section No. 1735, Vol. 18, No. 247, Pg. 87, May  
11, 1994 (19940511)

## ABSTRACT

PURPOSE: To execute a retrieval by a paragraph unit, and to execute an exact and high speed retrieval in accordance with a purpose by providing a means for discriminating the part corresponding to an item stored in a format storage means and retrieving its index.

CONSTITUTION: As for a document inputted by a document input means 1, in an index preparing means 2, paragraph-constituting information of the document classification concerned is obtained by referring to format data in a format storage means 5, and divided into plural paragraphs. These respective paragraphs are collected at every classification, and prepared as an index. The index generated by the index generating means 2 is stored by an index storage means 3. A fixed length paragraph index storage means 19 can manage index data by data of one table format, because data length of the paragraph is determined. In such a way, in the fixed length paragraph index storage means 19, a relational data base can be realized, and a high speed retrieval can be executed.



## 1

## 【特許請求の範囲】

【請求項1】 所定の書式を有する文書を検索する文書検索装置において、以下の要素を有する文書検索装置

- (a) 文書を記憶する文書記憶手段、
- (b) 文書の書式の構成を示す項目を記憶する書式記憶手段、
- (c) 文書記憶手段により記憶される文書から、書式記憶手段に記憶されている項目に対応する部分を識別して、その部分を項目に対応するインデックスとして記憶するインデックス記憶手段、
- (d) インデックス記憶手段により記憶されるインデックスを検索するインデックス検索手段。

【請求項2】 上記インデックス記憶手段は、項目に対応して、形式の異なるインデックスを有することを特徴とする請求項1記載の文書検索装置。

【請求項3】 所定の書式を有する文書を検索する文書検索装置において、以下の要素を有する文書検索装置

- (a) 文書を記憶する文書記憶手段、
- (b) 文書の書式の構成を示す項目を記憶する書式記憶手段、
- (c) 書式記憶手段により記憶される項目に対応する項目を書式辞書として記憶する書式辞書記憶手段、
- (d) 文書記憶手段により記憶されている文書と、書式記憶手段に記憶されている項目と、書式辞書記憶手段により記憶される書式辞書に基づき、インデックスを作成して記憶するインデックス記憶手段、
- (e) インデックス記憶手段により記憶されるインデックスを検索するインデックス検索手段。

【請求項4】 所定の書式を有する文書を検索する文書検索装置において、以下の要素を有する文書検索装置

- (a) 文書を記憶する文書記憶手段、
- (b) 文書の書式の構成を示す項目を記憶する書式記憶手段、
- (c) 文書記憶手段により記憶される文書から、書式記憶手段に記憶されている項目に対応する文書の部分を識別して、その部分を項目に対応するインデックスとして記憶するインデックス記憶手段、
- (d) インデックス記憶手段により記憶されるインデックスに対応する文書の部分を検索するための検索情報を記憶する検索情報記憶手段、
- (e) 検索情報記憶手段により記憶される検索情報を用いて、文書の部分を検索する部分検索手段。

## 【発明の詳細な説明】

## 【0001】

【産業上の利用分野】 この発明は蓄積された文書の中からキーワードを用いて、特に書式が予め決まっている文書を検索する文書の検索装置に関するものである。

## 【0002】

【従来の技術】 図22は第41回情報処理学会全国大会で発表されたフルテキストデータベース検索システム「検

## 2

蔵君」の文書検索装置を示すブロック構成図である。なお、図中実線は制御の流れを示し、点線はデータの流れを示している。図において、1は文書を入力するワープロなどの文書入力手段であり、入力された文書は磁気ディスク装置などの文書記憶手段13に格納される。8は検索要求として検索条件となるキーワードなどを入力する検索要求入力手段、14は入力されたキーワードに基づき文書記憶手段内の文書データを検索し、検索要求に適合する文書を限定する文書検索手段、10は使用者がどの文書を取り出すのかを指定する文書取り出し要求入力手段、11は文書取り出し要求に基づいて文書記憶手段13から対象文書を取り出す文書取り出し手段、12はどの文書が検索要求に適合したかという検索結果や取り出した文書データなどをディスプレイなどの出力装置に示す検索結果出力手段である。

【0003】 従来例の動作について、図22のブロック図を用いて説明する。文書そのものはワープロなどの文書入力手段1によって入力され、文書記憶手段13に記憶される。

【0004】 従来例の検索処理について説明する。文書内容との照合を行うためのキーワードなどの検索要求が検索要求入力手段8により入力される。入力された検索要求は、文書検索手段14へ送られ、文書検索手段14により文書記憶手段13内のすべての文書内容についてキーワードによる検索が行われる。検索の結果、検索要求に適合する文書が限定されるとそれらの文書情報は検索結果出力手段12へ送られ、検索結果出力手段12によりディスプレイなどに出力される。

【0005】 次に文書の取り出し処理について説明する。使用者は出力された検索結果を見て取り出したい文書を決め、文書取り出し要求手段10により文書番号や文書名などの文書取り出し要求を入力する。入力された文書取り出し要求は文書取り出し手段11に送られ、文書記憶手段13の中から文書取り出し要求に一致する文書のデータが文書取り出し手段11により取り出され、検索結果出力手段12へ送られる。検索結果出力手段12は送られてきた文書データをディスプレイなどに出力する。

【0006】 他の従来例について説明する。図23は、例えば、(株)平和情報センターの日本語文書情報系データベースシステム“Future/Happiness”のカタログに記載されているような従来のキーワードを用いた文書検索装置を示すブロック構成図である。なお、図中、実線は制御の流れを示し、点線はデータの流れを示している。

【0007】 図23を参照しながらこの従来例の構成を説明する。図において、1は文書を入力するワープロ等の文書入力手段であり、入力された文書は磁気ディスク装置等の文書記憶手段に格納される。231は格納された文書から自然言語処理機能によりキーワードを自動的

10

20

30

40

50

## 3

に抽出登録するキーワード登録手段であり、抽出されたキーワードは磁気ディスク装置等のキーワード記憶手段に格納される。8は検索要求となるキーワード等を入力するキーボードディスプレイ等の検索要求入力手段、14は入力されたキーワードに基づきキーワード記憶手段16を用いて検索要求に適合する文書を限定する文書検索手段、11はこの限定された文書の実体を、上記文書記憶手段13より取り出す文書取り出し手段であり、このようにして得られた検索結果はディスプレイ等の検索結果出力手段12により出力される。

【0008】次に動作について説明する。先ず文書そのものは文書入力手段1により入力され、文書記憶手段13に格納される。また同時に、キーワード登録手段231では入力された文書を文法に基づき、単語毎に分割する分かち書き処理を行った後、分割された単語を評価して、助詞等の不要な用語を除去することにより、検索時に必要となるキーワードを自動的に抽出し、キーワード記憶手段232に格納する。

【0009】次に、検索時について説明する。検索者は文書記憶手段13に記憶されている文書の中から自分の欲しい文書を検索する場合、検索要求入力手段8に検索条件としてキーワードを入力する。文書検索手段14は、キーワード記憶手段232に格納されているキーワードと、検索要求入力手段8で入力されたキーワードとのマッチングを行い、検索者の欲する文書を検索する。文書取り出し手段11では、文書検索手段14で検索された文書の実体を文書記憶手段13から取り出し、検索結果出力手段12であるディスプレイ等の表示装置に出力する。

【0010】この従来例の文書検索装置は以上のようにして、蓄積された大量の文書からキーワードをもとにして検索を行うものである。

## 【0011】

【発明が解決しようとする課題】従来の文書検索装置は以上のように構成されているので、キーワードなどの検索要求により、全文書の全範囲を検索する必要があり、検索に時間がかかるという問題点があった。また、入力されたキーワードをもとに全文書にわたって検索されるために、検索者が本来期待しない文書部分でのマッチングが行われ、不要な文書まで出力されてしまうという問題点があった。さらに、検索結果として文書を取り出す場合、文書単位でしか取り出せず、検索者は、その文書の中から自分の見たい必要なところを探さなければならなかった。

【0012】また、従来の文書検索装置は、以上のように構成されていたので、文書データに対して付加したキーワードを基に検索を行うことはできても、このキーワードが数および種別を限定されているために、キーワードで文書の内容を詳細に表現することは困難であり、文書の内容を示す詳細な事項を指定して文書の検索を行う

## 4

ことはできないという問題点があった。また、キーワードを基に検索を行う場合、上記に示したようにキーワードでは多様な意味を持つ文書の内容をすべて表現することはできないため、検索者が欲したものの他に、不要なものまで検索されてしまうという問題点があった。

【0013】この発明はこのような問題点を解消するためになされたものであり、高速に文書の検索ができ、文書中の必要な部分だけを取り出せる文書検索装置を得ることを目的とする。

## 10. 【0014】

【課題を解決するための手段】この発明の請求項1に係る文書検索装置は、検索すべき対象となる文書を記憶する文書記憶手段と、検索する文書の書式がどのような構成であるかを示す複数の項目を記憶する書式記憶手段と、文書記憶手段に記憶されている文書の内容を見て、書式記憶手段に記憶されている複数の項目に対応して文書の部分が、文書の中にあるかどうかを識別し、項目に対応する部分が合ったなら、その部分を項目に対応するインデクスとして記憶するインデクス記憶手段と、インデクス記憶手段により記憶されるインデクスの中に存在する検索すべき対象となる部分を検索するインデクス検索手段を備えたものである。

【0015】この発明の請求項2に係る文書検索装置は、請求項1のインデクス記憶手段において、項目の種別に対応して、形式の異なる複数種類のインデクスを有することができるものである。

【0016】この発明の請求項3に係る文書検索装置は、検索すべき対象となる文書を記憶する文書記憶手段と、検索する文書の書式がどのような構成であるかを示す複数の項目を記憶する書式記憶手段と、書式記憶手段により記憶される項目を特定の群として定める書式辞書を設けて、それを記憶する書式辞書記憶手段と、文書記憶手段に記憶されている文書の内容と、書式記憶手段により記憶されている書式辞書とに基づき、書式記憶手段の項目を書式辞書の項目と比べ、書式に項目に対応する部分が文書にあるかどうかを識別し、それにより、文書の部分を項目に対応するインデクスを作成して、記憶するインデクス記憶手段と、インデクス記憶手段により記憶されるインデクスの中に存在する検索すべき対象となる部分を検索するインデクス検索手段を備えたものである。

【0017】この発明の請求項4に係る文書検索装置は、検索すべき対象となる文書を記憶する文書記憶手段と、検索する文書の書式がどのような構成であるかを示す複数の項目を記憶する書式記憶手段と、文書記憶手段に記憶されている文書の内容を見て、書式記憶手段に記憶されている複数の項目に対応して文書の部分が、文書の中にあるかどうかを識別し、項目に対応する部分が合ったなら、その部分を項目に対応するインデクスとして記憶するインデクス記憶手段と、インデクス記憶手段に

## 5

より記憶される文書のある部分を検索するために、インデクスにおける文書の部分が、文書記憶手段に記憶されている文書のどの部分に存在するかという位置情報を記憶している検索情報記憶手段と、検索情報記憶手段により記憶される位置情報を用いて、文書記憶手段により記憶される文書の部分を検索する部分検索手段を備えたものである。

## 【0018】

【作用】この発明の請求項1の文書検索装置においては、文書が一定の書式に従って書かれている点に着目し、インデクス記憶手段はあらかじめ入力しておいた文書書式を用いて、文書データを論理的にまとまりのあるパラグラフに分割し、インデクスとして格納し、このインデクスを検索することにより、文書全体を検索の対象とせず、パラグラフ単位で検索ができる。また、検索時に、このようなインデクスを用いることにより、パラグラフを検索することになり、文書の内容に関する詳細な事項を指定して検索を行うことができる。

【0019】この発明の請求項2の文書検索装置においては、インデクス記憶手段は、インデクスの作成時、例えば、特許明細書における「書類名」、「発明の名称」などの構成を示す項目において、その項目に記される内容文章の長さが定まるパラグラフについては、固定長であるパラグラフとして、その他長さが定まらないものについては、可変長であるパラグラフとして、格納することができるようにインデクスの種類を分けるので、検索要求に対して、どちらかの適切なインデクスを選択してパラグラフの検索を行うことができる。

【0020】この発明の請求項3の文書検索装置においては、辞書を設けた、書式内に含まれない任意の項目を辞書に登録しておくことにより、書式内の項目と同様に扱い、それに従い、インデクスを検索できる。

【0021】この発明の請求項4の文書検索装置においては、検索情報記憶手段の記憶する位置情報に従い、部分検索手段により、検索要求に適合した文書に関して適合したパラグラフのみを検索結果として文書記憶手段から取り出すことができるので、必要な情報のみ得ることができる。

## 【0022】

## 【実施例】

実施例1. 図1は、実施例1におけるの文書検索装置の全体構成を示すブロック図である。従来例の図22と同一または相当部分には同一符号を用い、その説明は省略する。図1において、2は入力された文書を書式に基づいて、論理的な意味を持つ文章のかたまりであるパラグラフに分割して、パラグラフの種類毎にまとめたインデクスを作成するインデクス作成手段である。3はインデクス作成手段2により、パラグラフの種類ごとにまとめたインデクスを記憶する磁気ディスクなどのインデクス記憶手段である。4は、1つの文書がどのようなパラ

## 6

グラフで構成されているのかを示す書式を、文書の種類ごとに入力する書式入力手段、5は、その書式を記憶する磁気ディスク装置などの書式記憶手段である。6はインデクス作成手段2によって分割されたパラグラフが、文書記憶手段13に記憶されている元の文書においてどの部分だったのか、そして、インデクス記憶手段3内のどの場所に記憶されているのかというパラグラフの検索の目安となる位置情報を作成する検索情報作成手段、7は検索情報作成手段6により得られた情報を記憶する磁気ディスク装置などの検索情報記憶手段である。8は前記図22で示した検索要求入力手段8を改良したもので、検索要求を1つの文書全体に対してだけでなく、1つの文書内の各パラグラフに対しても出すことができるようになってい

る。9は検索要求手段8で入力された検索要求に対し、インデクス記憶手段3内のインデクスのパラグラフを検索し、1パラグラフの検索ごとに検索範囲の絞り込みを行うインデクス検索手段である。10は前記図22で示した文書取り出し要求入力手段10を改良したもので、文書取り出し要求を文書単位でなく、文書内のパラグラフ単位で指定できるようにしたものである。11は前記図22で示した文書取り出し手段11を改良したもので、文書取り出し要求入力手段10の指定により、文書取り出しを文書単位でなく、パラグラフ単位で行えるようにした部分検索手段15を含んでいる。12は前記図22で示した検索結果出力手段12を改良したもので、文書取り出し手段11により取り出してきたデータを、文書単位だけでなく、パラグラフ単位でも出力できるようにしている。なお、新たに備えられたインデクス作成手段2や検索情報記憶手段7などは、計算機システムを構成するプロセッサとその上で動作するソフトウェアによって実現されている。

【0023】次に動作について説明する。前述したように、本実施例は文書が一定の書式に従って書かれている点に着目したものであり、例えば、規格書・仕様書などの技術文書は図2に示すようにそれぞれ定まった書式を有している。文書の書式は予め書式入力手段4から入力され、書式記憶手段5に格納される。このとき、文書の書式は、書式の同じ文書を扱う場合には1種類でよいが、種別の異なる文書を扱う場合には文書の種別によって書式が異なるために、その種別の数だけ入力される。これらの文書の種別や構成を示す書式データは、文書をパラグラフごとに分割する際や検索要求入力時に、その種別が、決められた指定を用いて行われる。

【0024】一般に文書は、パラグラフ（ある意味に従い、かたまりとみなされる論理的な単位）の集まりからなる。このパラグラフには図3の例でいうと、「文書名」や「適用範囲」等のように、その内容文章の長さが何文字以内と定まっているもの（固定長パラグラフ）と、「一般要求事項」のようにその長さが何文字以内と

10

20

30

40

50

## 7

る。従って、書式として入力される書式データは、図4に示すように、文書がどのようなパラグラフから構成されているかを表す文書構造4 1、パラグラフの内容を表す書式項目名4 2、書式項目名に対して付けられた書式項目番号4 3、パラグラフが固定長か可変長かを表す固定長／可変長フラグ4 4、書式項目名に対して付けられるパラグラフID 4 5などにより構成されることになり、パラグラフ構成情報として対応するフラグで表現される。

【0025】文書そのものは従来例通り文書入力手段1によって入力される。入力された文書は、インデックス作成手段2において、書式記憶手段5内の書式データの参照により該当する文書種別のパラグラフ構成情報が得られ、複数のパラグラフへと分割される。更に、分割された各パラグラフは、図5のようにパラグラフの種別ごとにまとめられ、インデックスとして作成される。そして、インデックス作成手段2で作成されたインデックスは、インデックス記憶手段3により格納される。

【0026】次に、本実施例でのインデックスの詳細について説明する。図6は固定長パラグラフのインデックスを記憶する固定長パラグラフインデックス記憶手段の構造を示したものである。この固定長パラグラフインデックス記憶手段については、パラグラフのデータ長が定まるため、インデックスデータを1つの表形式のデータで管理することができる。このため固定長パラグラフインデックス記憶手段では、リレーショナルデータベースの実現が図れ、検索時に高速な検索を行うことができる。

【0027】図7は可変長パラグラフのインデックスを記憶する可変長パラグラフインデックス記憶手段の構造を示したものである。この可変長パラグラフインデックス記憶手段については、パラグラフのデータ長が定まらないため、各パラグラフのインデックスごとにデータを一つの表形式で管理する。

【0028】インデックス作成手段2により作成されるインデックスにおいてのパラグラフ分割の前の元文書に関する情報（文書名、作成者など）と、インデックス記憶手段3により記憶されるパラグラフ分割後に関する情報（文書中の位置情報、格納場所など）はリンクされ、図8に示すように、位置情報として検索情報作成手段6により作成され、検索情報記憶手段7に記憶される。

【0029】次に本実施例の検索処理について説明する。従来の検索要求入力とは文書記憶手段13に登録されている全文書に対し、それぞれの文書の内容すべてを検索対象として、キーワードを入力することで検索が行われていた。本実施例では、図9に示すように文書の種類を選択し、書式記憶手段5内の書式データのパラグラフ構成情報を参照することにより、各パラグラフを検索対象としてキーワードを入力する。入力されたキーワードは、文書の種別の情報や検索対象のパラグラフの種別の情報と一緒にインデックス検索手段9へと送られる。イン

## 8

デックス検索手段9では、検索要求手段8から送られてきた検索要求に基づいて、インデックス記憶手段3に記憶されたパラグラフをキーワードに基づき検索する。その結果、検索要求に適合する文書が限定される。

【0030】なお、定型な文書書式をもつ文書においては、文書構造、書式項目名、書式項目番号が図4に示したように書式データのパラグラフ構成情報として定まっているが、実際には図10に示すように、書き手が自分の使いやすいうように項目名を変え、書式を変えて文書を作成するケースがある。この様な場合には、同じ書式であっても、違う書式データを作成することになってしまう。本実施例の文書検索装置は、書き手によって文書書式が違う文書に対しても、インデックス作成手段2において、書式辞書及び書式辞書記憶手段を設けたので、その違いに対応することができる。書式辞書とは、定型な文書書式に対して、書き手が自分の作り易いように文書書式を変更して作成した文書からでも、インデックスの作成を可能にするために設けられたものである。具体的には図11に示すように、予め定まっている書式項目名に対して、書き手が変更して書き得る可能性のある項目名を辞書としたものであり、図4のパラグラフ構成情報のパラグラフID 4 5を書式項目名に対応して蓄積している。

【0031】図12はインデックス作成手段2の構成を示すブロック図である。16は書式記憶手段5に記憶されている文書書式の各項目に対応した書式項目名をあらかじめ辞書として入力する、ワープロ等の書式辞書入力手段、17は入力された書式辞書を記憶する磁気ディスク装置等の書式辞書記憶手段、18は文書記憶手段13に蓄積されている文書データの中からその中に書かれている書式項目名を抽出し、抽出された書式項目名と書式辞書記憶手段17に記憶されている書式辞書の項目名とを比較することによって、分割されたパラグラフと文書書式との対応付けを行うインデックス作成手段である。なお、インデックス作成手段2は、計算機システムを構成するプロセッサとその上で動作するソフトウェアによって実現されている。

【0032】インデックス作成の動作について図12を用いて説明する。インデックス作成手段2において、文書入力手段1から文書が入力され文書記憶手段13に格納される度に、書式記憶手段5に記憶されている文書の書式データを用いて、文書記憶手段13の文書データをパラグラフに分割し、インデックスを作成する。

【0033】上記のインデックス作成手段2は、図13に示すような一連の処理を行う。以下に、それぞれの処理について図14の文書例をもとに説明を行う。

【0034】(a) 書式項目名抽出処理131；文書記憶手段13に蓄積されている文書データの中から、新たに登録された文書データを先頭行から読みだし、書式項目番号（図4の43参照）を検出した場合は、その次に

書かれている書式項目名を抽出する。図14の文書では、先ず先頭行として「1. 文書名称」が読み込まれ書式項目番号「1.」が検出され、続く「文書名称」が書式項目名として抽出される。

【0035】(b) パラグラフID獲得処理132;

(a) で抽出された書式項目名をキーとして書式辞書の検索を行い、マッチングしたデータのパラグラフIDを得る。図11に示した書式辞書を例にとると、図14の文書では、書式項目名「文書名称」はパラグラフIDが「1」に対応することが分かる。

【0036】(c) パラグラフ抽出処理133; (a) で書式項目名が検出された次の行から、また書式項目名が検出されるまでデータを読み込む。この際読み込まれたデータがパラグラフである。図14の文書では、次の書式項目番号「適用する範囲」が検出されるまでのデータ、すなわち、2行目の「△△△仕様書」がパラグラフデータとして抽出される。

【0037】(d) インデックス作成・登録処理134; (b) で得られたパラグラフIDおよび、(c) で抽出されたパラグラフデータを図15に示すように対応付けてインデクスデータとしてインデクス記憶手段に出力する。なお、本実施例では、検索の際に高速化を図るため、先に述べたように固定長パラグラフと可変長パラグラフの場合でインデクスの記憶手段を分けている。インデクスの登録の際には、書式記憶手段5に記憶される書式データのパラグラフ構成情報(文書構造、書式項目名、書式項目番号、パラグラフID、固定長/可変長フラグ)を参照し、(b) で得られたパラグラフIDをキーとして、該当する書式項目の固定長/可変長フラグをチェックし、(c) で得られたパラグラフを固定長パラグラフインデクス記憶手段に登録するか、可変長パラグラフインデクス記憶手段に登録するかを判定している。

【0038】(e) 継続行の判定135; (a) ~

(d) の処理の後に、文書内に継続行がある場合は、継続行がなくなるまで繰り返して処理する。

【0039】次に検索時について説明する。例えば、図16(a)、(b)に示すようなインデクスデータが、上述したインデクス作成手段2によって作成され、それぞれ固定長パラグラフインデクス、可変長パラグラフインデクスに格納される。また図16(c)に示すような書式データが書式記憶手段5に格納される。このような場合に、検索者が「書式項目『適用範囲』の中に『文書検索装置』というワードがあるもの」というように文書の内容を指定して検索を行いたいときは、検索要求入力手段8において、図17(a)に示すように書式記憶手段5に記憶されている書式を画面上に表示し、検索者はシステムが表示したこの書式の書式項目毎に、文書の内容を表すワードをキーワードとして入力することにより検索が行える。上記の例では、書式項目『適用範囲』に対して、キーワード『文書検索装置』を入力することに

なる。このようにして入力された、文書の内容を指定した検索条件は、図17(b)に示すように、書式データのパラグラフIDとキーワードを組として、文書名が検索結果として(c)のように検索結果手段12に出力される。なお、この際、同時に、書式データのうち指定された書式項目の固定長/可変長フラグも同時に出力される。

【0040】ここで、文書取り出し手段11は、インデクス記憶手段3に記憶されたインデクスデータを用いて、検索要求入力手段8より入力された検索条件をもとに検索を行う。図17の例では、先ず、検索条件の中でキーワードが設定された書式項目の固定長/可変長フラグを判定する。図17の例では検索条件の中で、固定長パラグラフを表すフラグが示されているので、固定長パラグラフインデクスを用いて検索を行う。具体的には、固定長パラグラフインデクスの中で、パラグラフIDが「2」の列に存在するパラグラフの中に、キーワードとして入力したワードが存在するかどうかを調べ、存在した場合には検索条件に適合したものとみなす。この場合、図17(c)に示すように文書Bが検索結果として得られるが、さらに、文書取り出し手段11により、文書記憶手段13から文書の内容が取り出され、検索結果出力手段8であるディスプレイ等の表示装置に出力される。

【0041】また、上記では固定長パラグラフに対応した書式項目にキーワードが設定された場合について説明したが、可変長パラグラフに対応した書式項目にキーワードを設定することももちろん可能である。図18は上記のような検索の具体例を示したものであり、可変長パラグラフに対応した書式項目『一般要求事項』の中に『マルチメディア』というワードがある文書を検索する場合を示している。検索要求の入力は、図18(a)に示すように、書式項目が固定長パラグラフに対応したものと同様に、書式項目『一般要求事項』に対して、キーワード『マルチメディア』を入力するだけでよい。従って、検索者はキーワードを入力しようとしている書式項目が固定長か可変長かについては意識する必要がない。図18(b)は検索条件を示している。この場合、固定長/可変長フラグが可変長に設定される。

【0042】ここで、文書取り出し手段11は、インデクス記憶手段に記憶されたインデクスデータを用いて、検索要求入力手段8より入力された検索条件をもとに検索を行う。この場合、可変長パラグラフを表すフラグが示されているので、可変長パラグラフインデクスを用いて検索を行う。具体的には、可変長パラグラフインデクスデータの中でパラグラフIDが「4」のパラグラフデータの中にキーワードとして入力したワードが存在するかどうかを調べ、存在した場合には検索条件に適合したものとみなす。この場合、図18(c)に示すように文書A、Bが検索結果として得られる。以上の検索により

検索結果として得られた文書は、文書取り出し手段11により文書記憶手段13より取り出され、検索結果出力手段12であるディスプレイ等の表示装置に出力される。

【0043】本実施例での部分検索手段15について説明する。検索者は出力された検索結果を見て、自分の取り出したい文書を決め、文書取り出し要求手段10により、文書取り出し要求を入力する。このとき、従来方式では文書単位でしか内容の取り出しが指定できなかったのに対し、本実施例では図19に示すように、文書の10 パラグラフまで指定して取り出し要求が入力できる。入力された文書取り出し要求は文書取り出し手段11に送られ、要求に適合した文書が取り出される。このとき、文書取り出し手段11は図8に示す検索情報記憶手段7内の位置情報を参照して、部分検索手段15により、文書記憶手段13内の文書の内容をパラグラフ単位で取り出すことができる。部分検索手段15により得られたパラグラフは、検索結果出力手段12により出力される。

【0044】実施例2. 上記実施例では、一つの書式項目に対してキーワードを入力する場合を説明したが、複数の書式項目に対してキーワードを入力し、検索すること10 ももちろん可能である。図20は上記のような検索の具体例を示したものであり、書式項目『適応範囲』の中に『装置』というワードがあり、しかも、書式項目『一般要求事項』の中に『マルチメディア』というワードがある文書を検索する場合を示している。この場合は、固定長パラグラフインデックス記憶手段19と可変長パラグラフインデックス記憶手段20を用いることにより検索が行なわれる。

【0045】実施例3. 実施例1の検索要求入力手段830 において、複数のパラグラフを検索対象としてキーワードの入力が行われた場合、インデックス検索手段9は、検索情報記憶手段7を参照して、1パラグラフ検索ごとにインデックス内における検索範囲の絞り込みを行いながら検索を行うことが可能である。図21のようにインデックスA、B、Cに対して結合条件ANDでそれぞれのパラグラフに対してあるキーワードが入力された場合、インデックス検索手段9はまずインデックスAの検索を行う。インデックスAに対する検索の結果、文書1・文書3・文書4が条件に適合したとする。次のインデックスBの検索を行う際に、インデックス検索手段9は検索情報記憶手段7を参照して、インデックスB内の文書1・文書3・文書4の部分のみを検索対象として取り出し、検索を行う。以降はこの繰り返しにより検索範囲の絞り込みが行われる。なお、検索要求の結合条件がORの場合は、条件に適合した文書以外の文書のパラグラフを次の検索対象とすることにより、絞り込みが行われる。そして、インデックス検索手段9により限定された文書のパラグラフは、検索結果出力手段12によりディスプレイなどに出力される。このように、検索要求が複数のパラグラフに対し

である場合、パラグラフごとに検索が行われる際に、前のパラグラフの検索結果から検索要求に適合する文書が限定されるので、検索情報記憶手段を参照することで、その次のパラグラフ検索では前回の検索時より範囲を限定した検索を行うことができる。

【0046】実施例4. 実施例1では、インデックス記憶手段において、文書の分割されたパラグラフを有するとしたが、これは、文書記憶手段に対して文書検索を行うよりも、インデックス記憶手段に対して検索を行うほうが、より高速に検索できるためである。しかし、この方法では、記憶容量の問題があると思われる場合には、インデックス記憶手段においてパラグラフを有さず、そのかわり、パラグラフの存在する文書記憶手段内の文書の位置を情報としてインデックス記憶手段で有することで、インデックス記憶手段のパラグラフの位置情報を元に、常に文書記憶手段に対して文書を検索する方法を取ってもよい(ただし、この方法では、時間的問題がある)。また、インデックス記憶手段では、文書記憶手段に記憶されている文書のパラグラフそのものと全く同じものを記憶していてもよいし、あるいは、文書記憶手段に記憶されている文書のパラグラフにおいて、野線情報や特殊な文字コード情報などの不要な情報を除いた部分のみを記憶しているものであってもよい。

【0047】実施例5. 実施例1において、インデックス記憶手段では、検索の高速化を狙い、固定長パラグラフインデックス記憶手段と可変長パラグラフインデックス記憶手段とを分別して設けたが、特にこの2種類に分別する必要はなく、例えば、英文と和文の区別と言うように、他の種別による記憶手段を、インデックス記憶手段において、設けるものであってもよい。また、記憶手段の分別は、2種類以上であってもよい。

【0048】実施例6. 実施例1において、書式辞書記憶手段は、インデックス作成手段が有するとして説明したが、特にインデックス作成手段が有さずに、独立した機能として存在してもよい。書式辞書入力手段においても、同様である。

【0049】実施例7. 実施例1において、書式記憶手段に記憶される書式データのパラグラフ構成情報は、図4に示すように、数字のフラグを使用しているが、例えば、文字や記号を使用するものでもよい。また、書式記憶手段に記憶される書式データの構成は、図4のような階層構造以外に、表形式などの形式でも構わない。

【0050】実施例8. 実施例1において、部分検索手段は、文書取り出し手段に含まれるとしたが、独立した機能として存在してもよい。また、部分検索手段の実施方法は図19に示したが、検索対象とするパラグラフの選択の方法に関しては、特に限定はしない。

【0051】

【発明の効果】以上のように請求項1の発明による文書50 検索装置では、文書を記憶する文書記憶手段と、文書の



13

書式の構成を示す項目を記憶する書式記憶手段と、文書記憶手段により記憶される文書から、書式記憶手段に記憶されている項目に対応する部分を識別して、その部分を項目に対応するインデクスとして記憶するインデクス記憶手段と、インデクス記憶手段により記憶されるインデクスを検索するインデクス検索手段とを設けたことにより、文書の検索時にパラグラフ単位で検索ができるので、検索者の目的に応じて、適格で高速な検索ができる。

【0052】以上のように請求項2の発明による文書検索装置では、請求項1のインデクス記憶手段が、項目に対応して、形式の異なるインデクスを有することにより、検索対象のパラグラフを固定長と可変長に分別して管理できるので、検索者の必要な文書のみを高速に検索することが可能である。

【0053】以上のように請求項3の発明による文書検索装置では、文書を記憶する文書記憶手段と、文書の書式の構成を示す項目を記憶する書式記憶手段と、書式記憶手段により記憶される項目に対応する項目を書式辞書として記憶する書式辞書記憶手段と、文書記憶手段により記憶されている文書と、書式記憶手段に記憶されている項目と、書式辞書記憶手段により記憶される書式辞書に基づき、インデクスを作成して記憶するインデクス記憶手段と、インデクス記憶手段により記憶されるインデクスを検索するインデクス検索手段とを設けたことにより、書式の類似している文書を同一書式とみなして検索できるので、検索者の目的に応じた文書を適格に検索できる。

【0054】以上のように請求項4の発明による文書検索装置では、文書を記憶する文書記憶手段と、文書の書式の構成を示す項目を記憶する書式記憶手段と、文書記憶手段により記憶される文書から、書式記憶手段に記憶されている項目に対応する文書の部分を識別して、その部分を項目に対応するインデクスとして記憶するインデクス記憶手段と、インデクス記憶手段により記憶されるインデクスに対応する文書の部分を検索するための検索情報を記憶する検索情報記憶手段と、検索情報記憶手段により記憶される検索情報を用いて、文書の部分を検索する部分検索手段とを設けたことにより、検索対象の文書を、文書の内容全体としてだけでなく、そのないようの一部分であるパラグラフも検索できるので、検索者の必要な情報のみを高速に得ることができる。

【図面の簡単な説明】

【図1】本発明の実施例1における文書検索装置の全体構成を示すブロック図である。

【図2】本発明の実施例1で用いられる一定の書式に従って書かれている文書の一例を示す図である。

【図3】本発明の実施例1の書式と文書データにおけるパラグラフの対応を示す図である。

【図4】本発明の実施例1の書式データにおけるパラグ

14

ラフ構成情報の例を示す図である。

【図5】本発明の実施例1のインデクス作成手段により分割された文書のパラグラフが、インデクス記憶手段により格納される例を示す図である。

【図6】本発明の実施例1の固定長パラグラフインデクス記憶手段の構成を示す図である。

【図7】本発明の実施例1の可変長パラグラフインデクス記憶手段の構成を示す図である。

【図8】本発明の実施例1のインデクス記憶手段と検索情報記憶手段の関係を示す図である。

【図9】本発明の実施例1の検索処理の流れを示す図である。

【図10】本発明の実施例1の文書書式の一例を示す図である。

【図11】本発明の実施例1の文書書式と書式辞書の対応を示す図である。

【図12】本発明の実施例1のインデクス作成手段の構成を示す図である。

【図13】本発明の実施例1のインデクス作成手段の処理の流れを示す図である。

【図14】本発明の実施例1の文書例を示す図である。

【図15】本発明の実施例1の文書データと書式データにおけるパラグラフIDの対応を示す図である。

【図16】本発明の実施例1の固定長パラグラフインデクスデータと可変長パラグラフインデクスデータにおけるパラグラフIDの書式データへの対応を示す図である。

【図17】本発明の実施例1の固定長パラグラフに対応した検索方法の一例を示す図である。

【図18】本発明の実施例1の可変長パラグラフに対応した検索方法の一例を示す図である。

【図19】本発明の実施例1の部分検索手段によるパラグラフの検索方法を示す図である。

【図20】本発明の実施例2の検索方法を示す図である。

【図21】本発明の実施例3の検索においてインデクスのパラグラフの絞り込み方法を示す図である。

【図22】従来例の文書検索装置の全体構成を示すブロック図である。

【図23】従来例の文書検索装置の全体構成を示すブロック図である。

【符号の説明】

- 1 文書入力手段
- 2 インデクス作成手段
- 3 インデクス記憶手段
- 4 書式入力手段
- 5 書式記憶手段
- 6 検索情報作成手段
- 7 検索情報記憶手段
- 8 検索要求入力手段

15

16

9 インデクス検索手段

10 文書取り出し要求入力手段

11 文書取り出し手段

12 検索結果出力手段

13 文書記憶手段

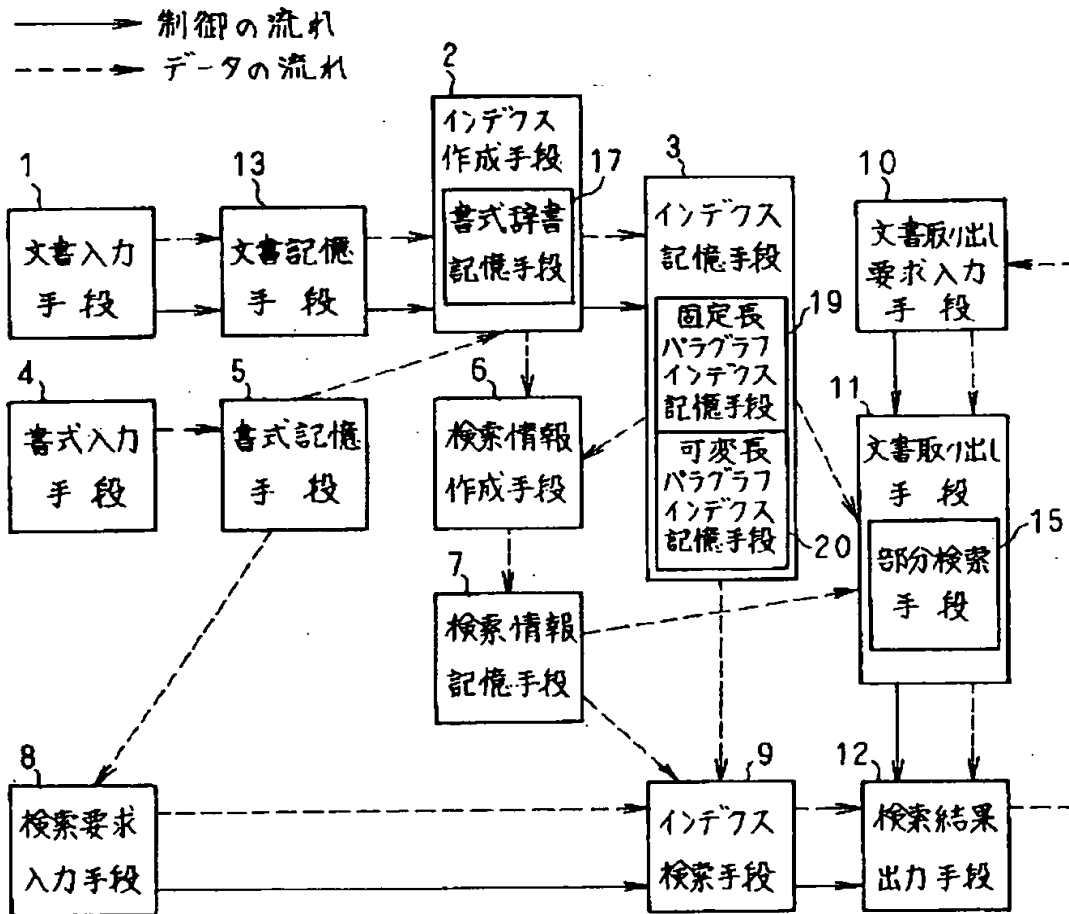
15 部分検索手段

17 書式辞書記憶手段

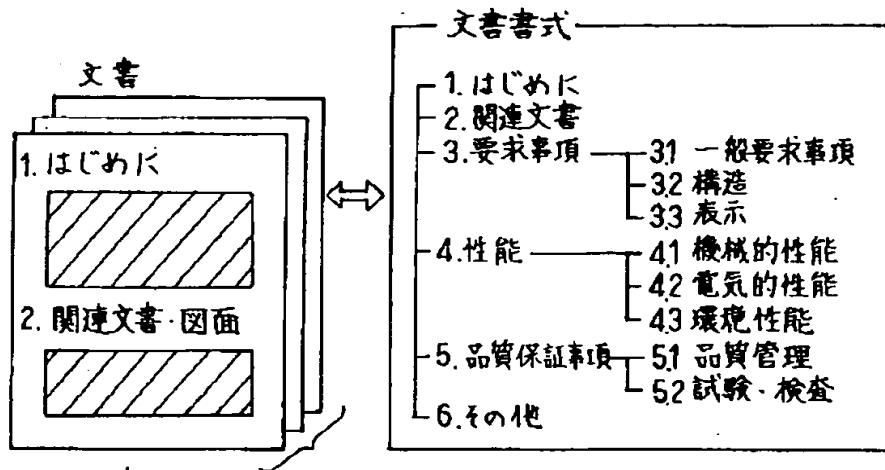
19 固定長パラグラフインデクス記憶手段

20 可変パラグラフインデクス記憶手段

【図1】

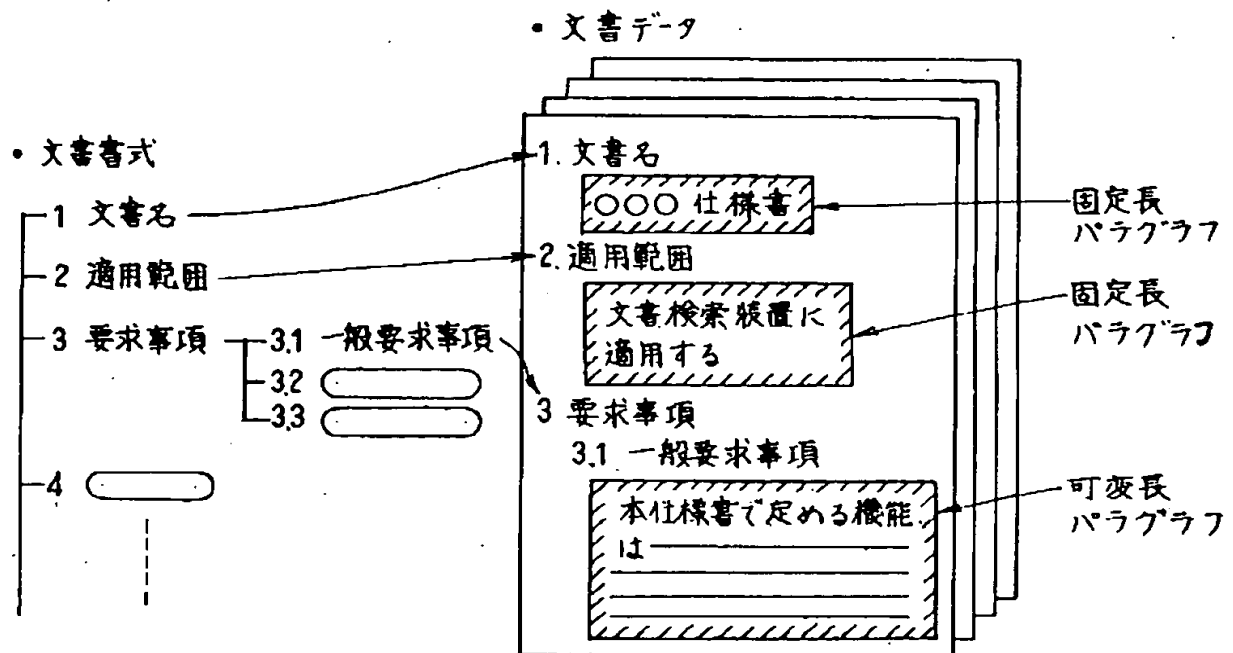


【図2】

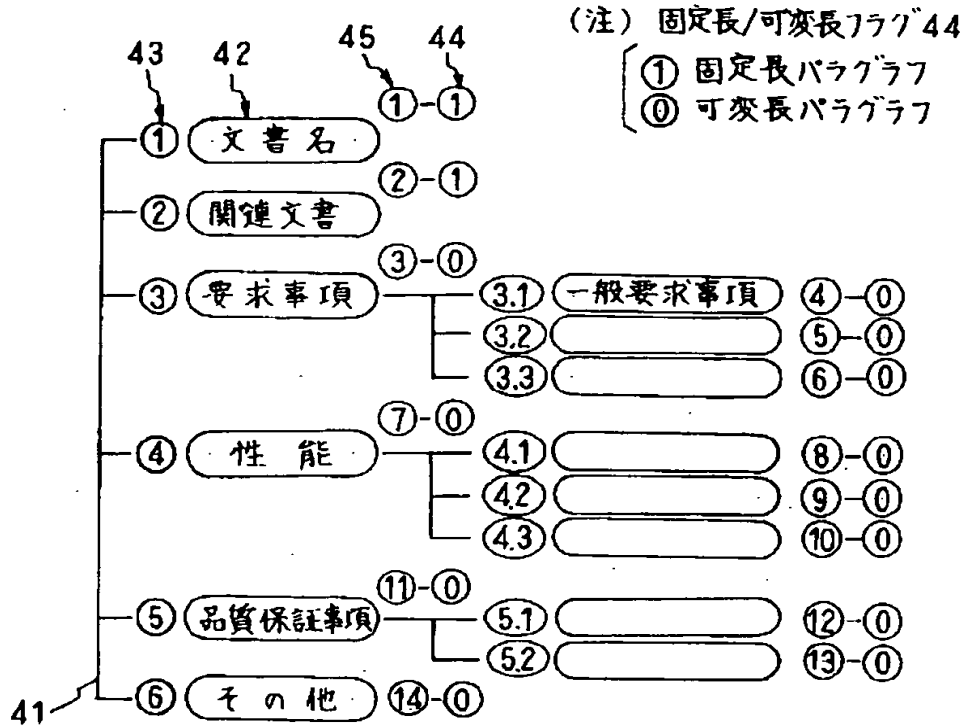


(例. ○○報告, ○○規格)

【図3】

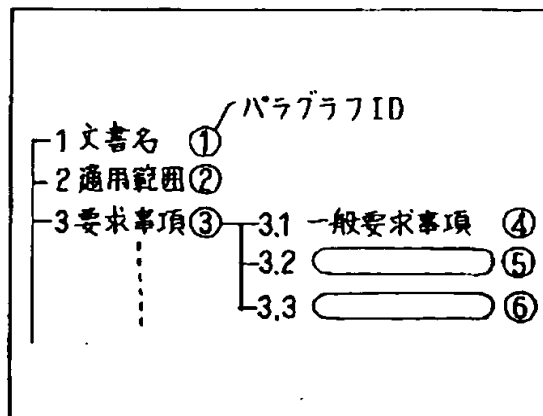


【図4】



【図11】

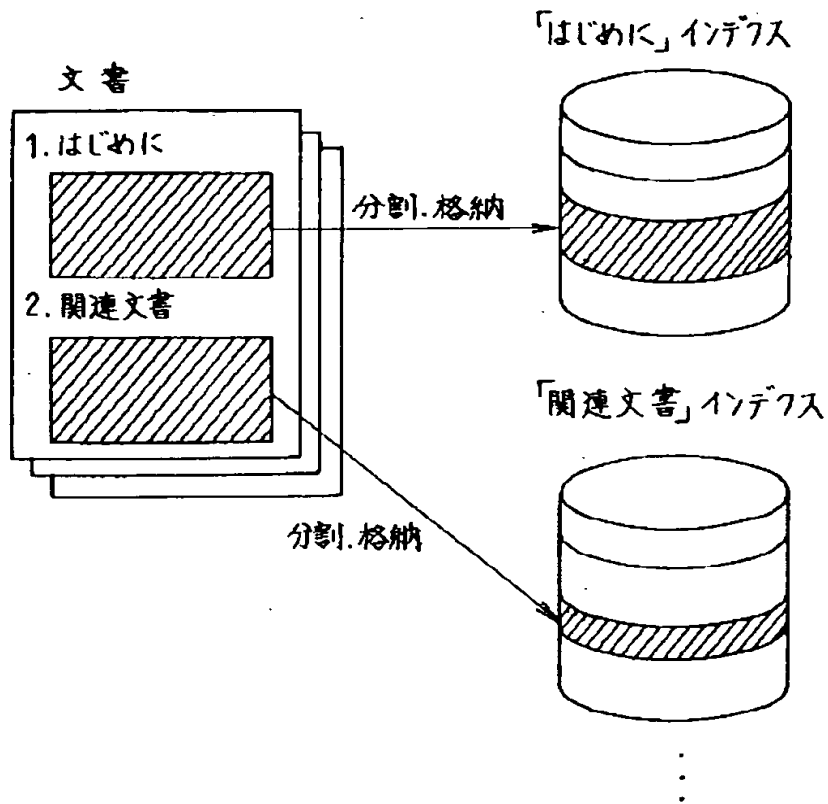
## 文書書式



## 書式辞書

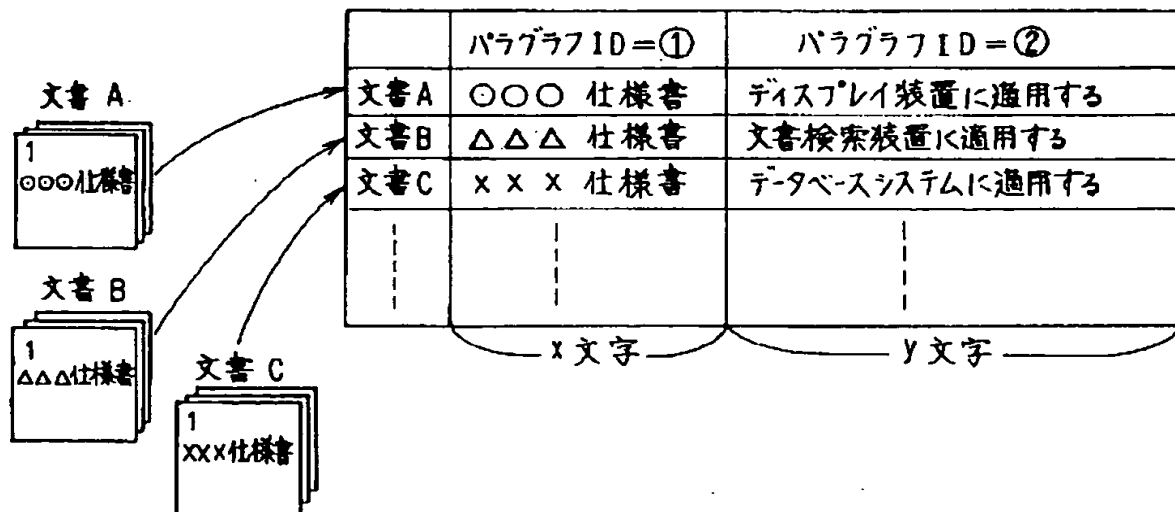
書式項目名	パラグラフID
文書名	①
文書名称	①
仕様書名	①
...	...
適用範囲	②
適用する範囲	②
...	...

【図 5】

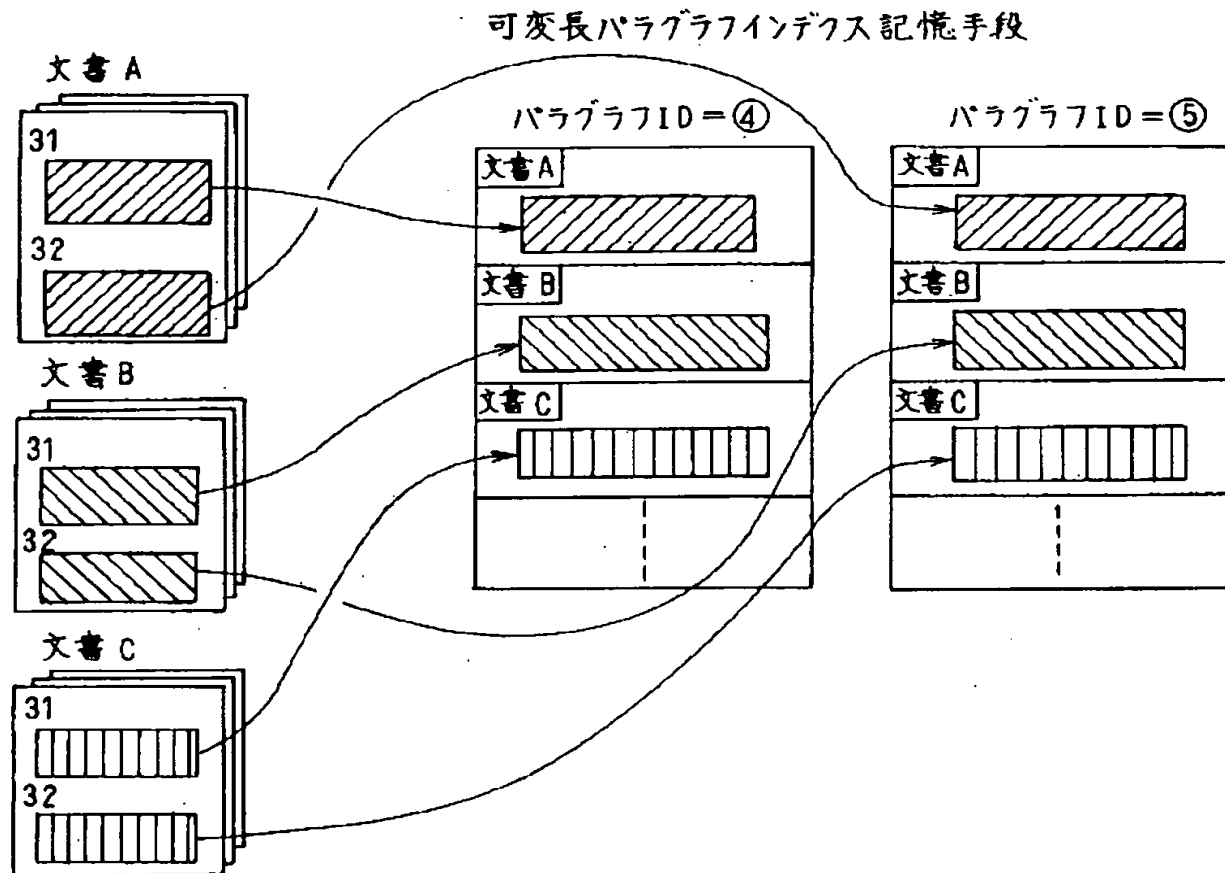


【図 6】

固定長パラグラフインデックス記憶手段



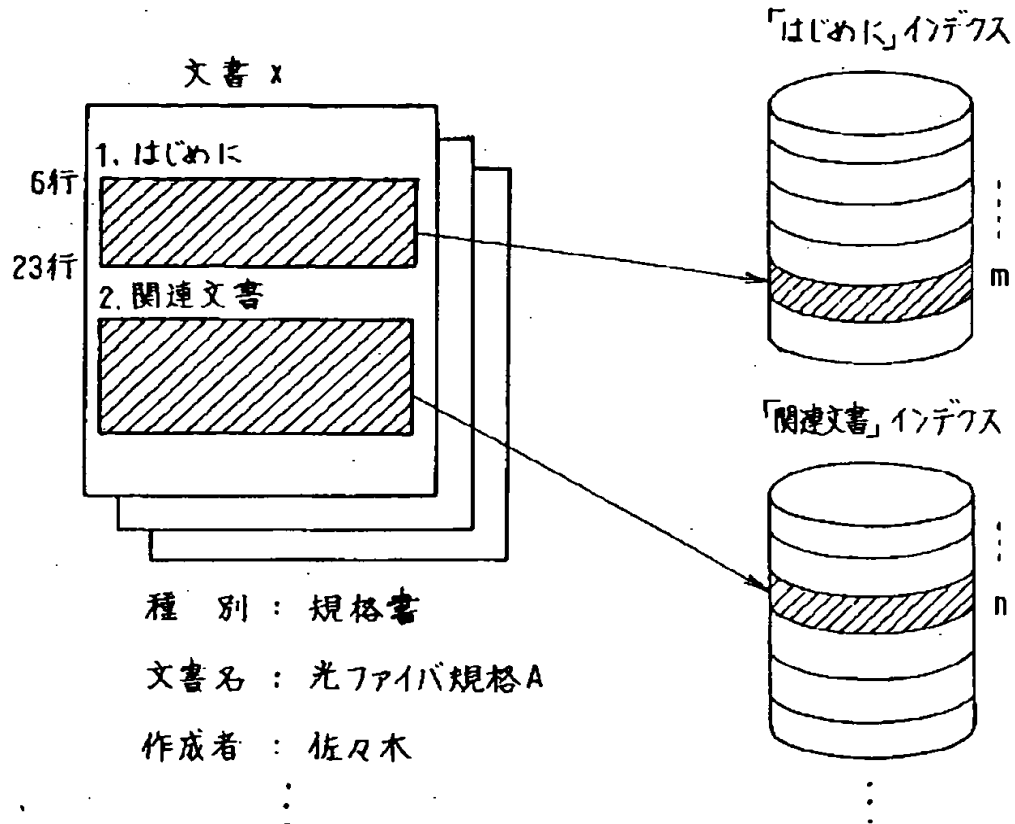
【図7】



【図14】

行数	
1	1 文書名称
2	△△△ 仕様書
3	2 適用する範囲
4	文書検索装置に
5	適用する
6	3 一般要求事項
7	本仕様書で定める機能
8	は _____
	_____
	_____

【図8】



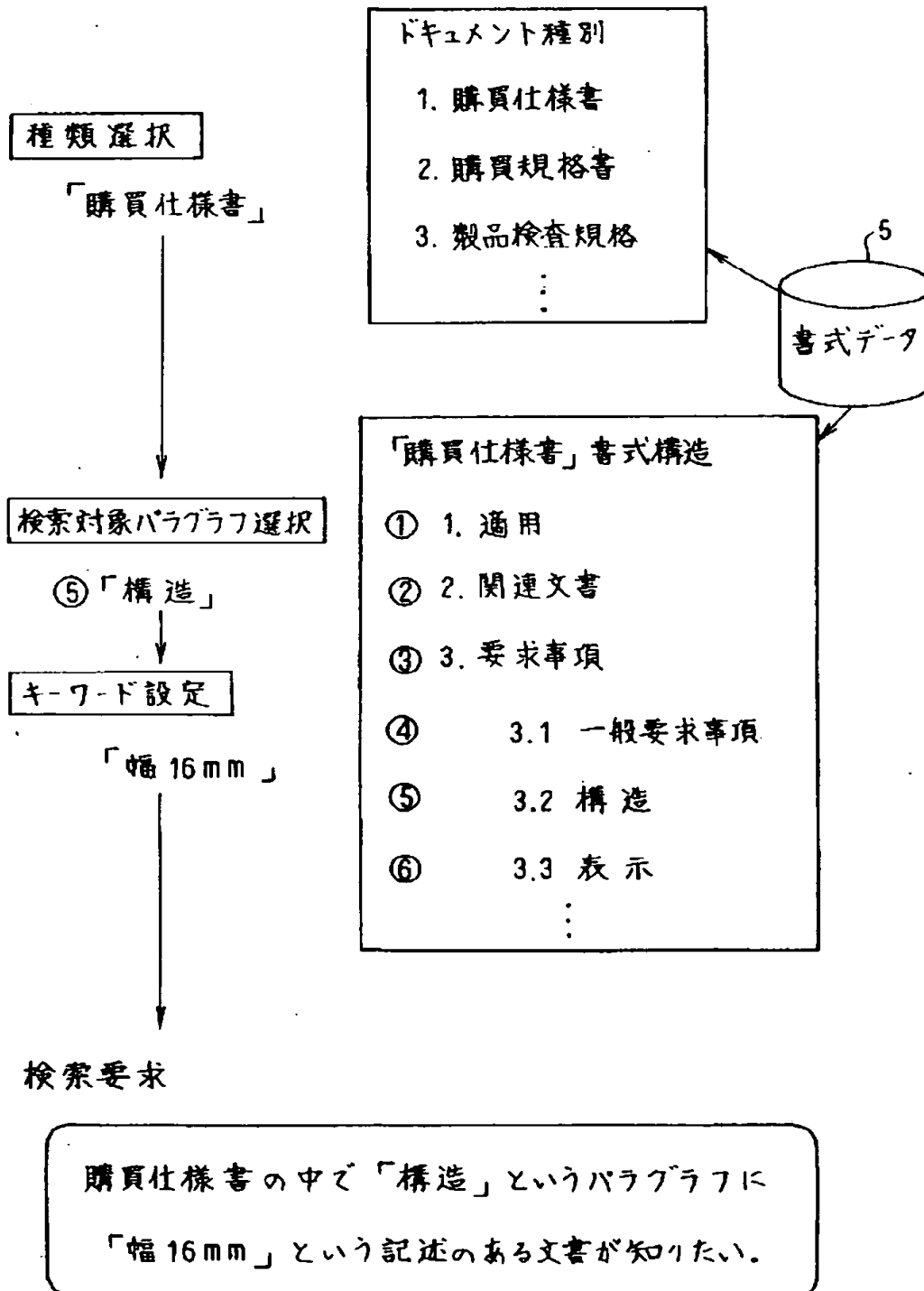
↓ 6 検索情報作成手段

「はじめに」位置情報

文書	文書内位置	格納位置
文書 1	1ページ 8~20行	5 ブロック
文書 2	2ページ 1~15行	12 ブロック
⋮	⋮	⋮
文書 X	1ページ 6~23行	m ブロック
⋮	⋮	⋮

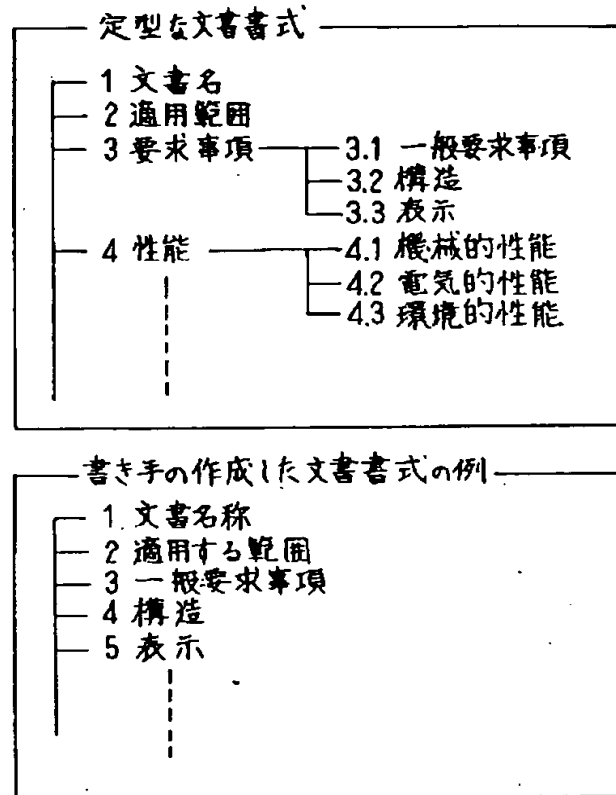
7  
検索情報  
記憶手段

【図9】

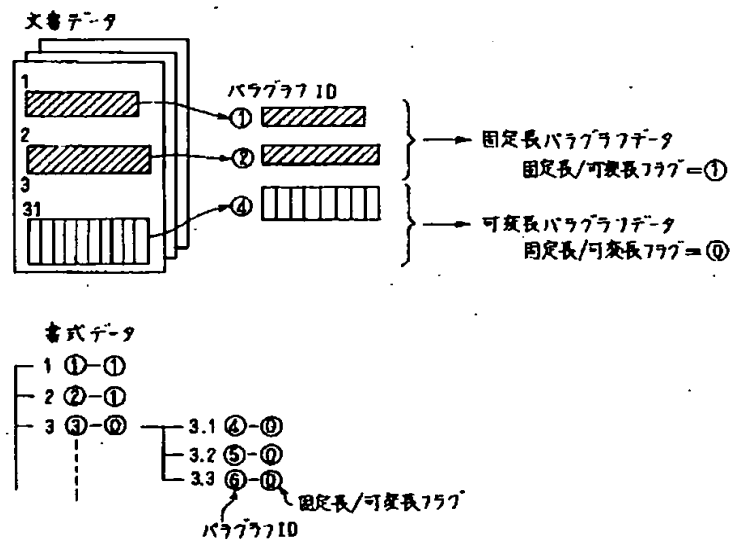




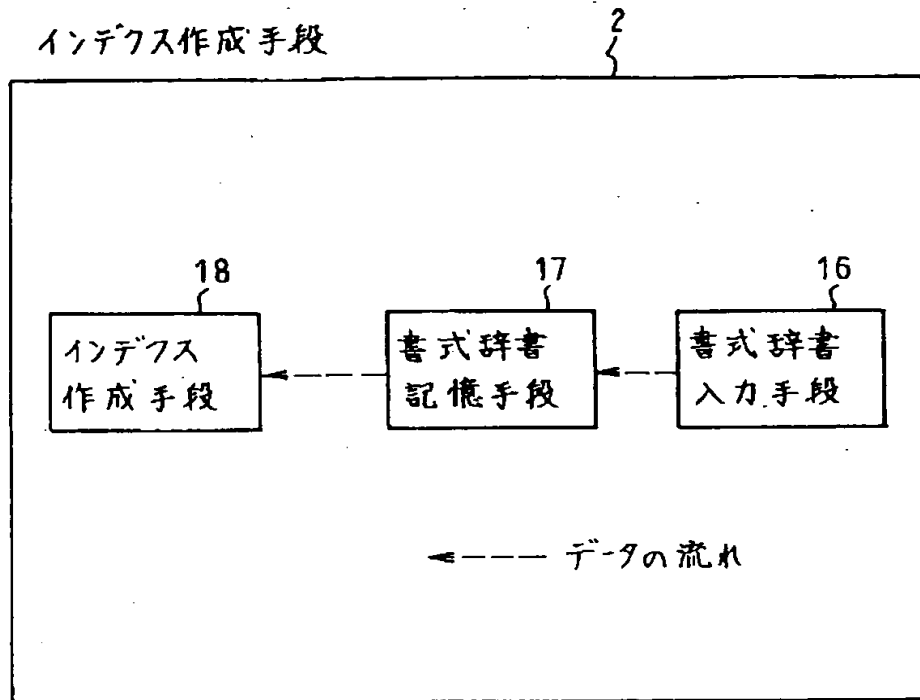
【図10】



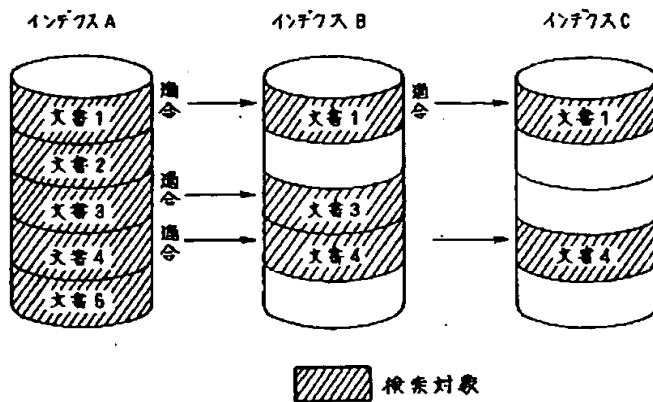
【図15】



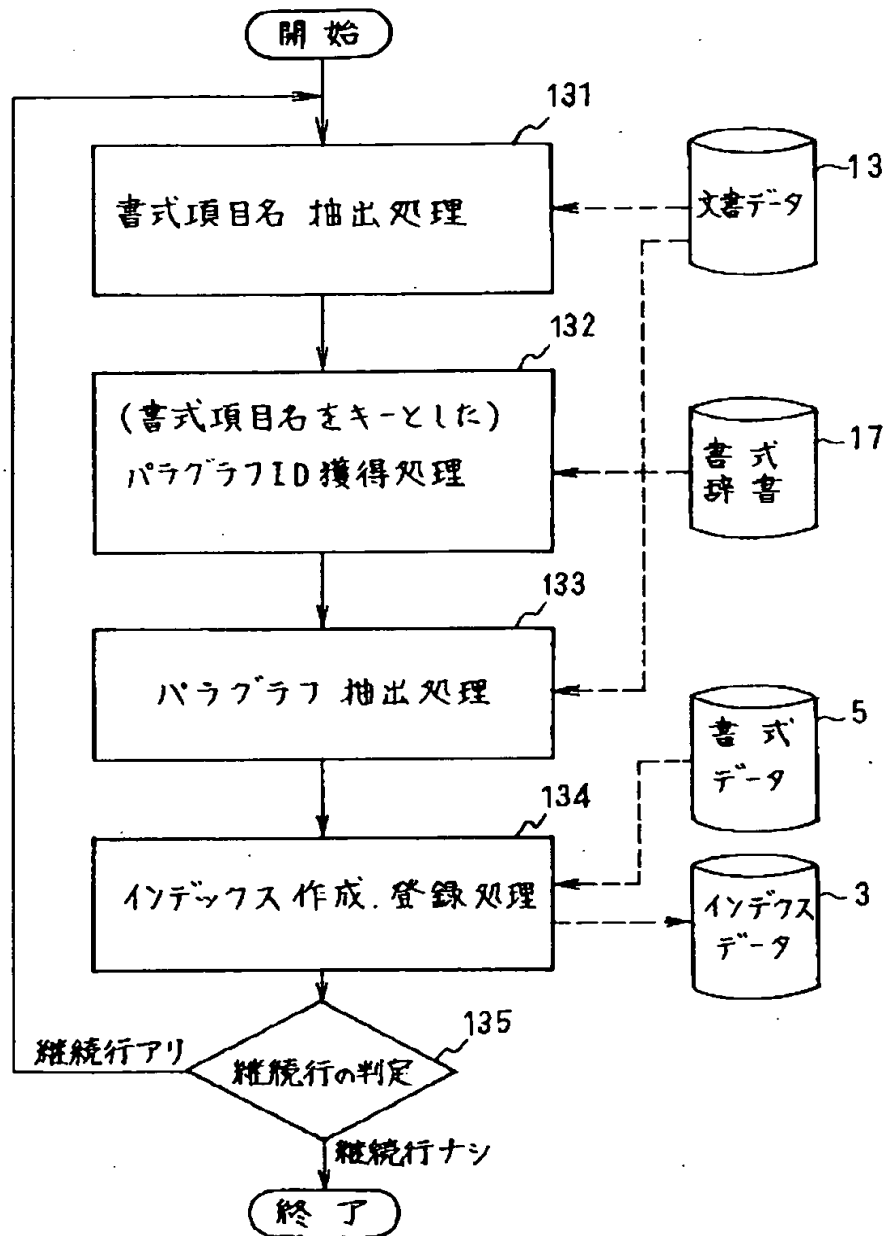
【図12】



【図21】



【図13】



【図16】

## (a) 固定長パラグラフインデクスデータ

パラグラフID=①

パラグラフID=②

文書A	〇〇〇仕様書	ディスプレイ装置に適用する
文書B	△△△仕様書	文書検索装置に適用する
文書C	×××仕様書	データベースシステムに適用する

## (b) 可変長パラグラフインデクスデータ

パラグラフID=④

文書A	本仕様書では画像データを表示する-----
文書B	本仕様書ではマルチメディア文書を-----
文書C	本仕様書ではマルチメディアデータを-----

パラグラフデータID=⑤

文書A	外形寸法は、たて950mm よこ300mm-----
文書B	サーバマシンのデータは 1GBのディスク-----
文書C	データとしては画像、図 形、文字-----

## (c) 書式データ

- 1 文書名 ①-①
- 2 適用範囲 ②-①
- 3 要求事項 ③-⑦

- 3.1 一般要求項目
- 3.2 構造
- 3.3 表示

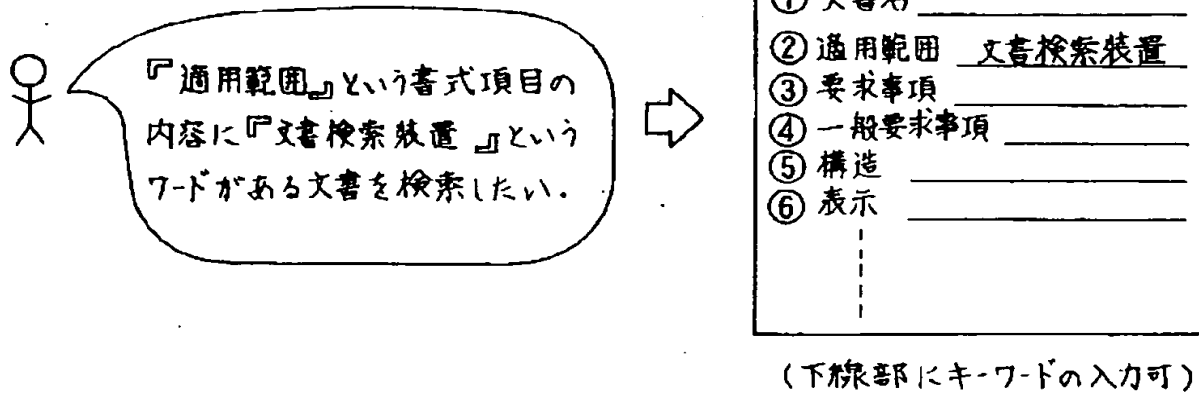
- ④-⑦
- ⑤-⑦
- ⑥-⑦

パラグラフID

固定長/可変長フラグ

【図17】

## (a) 検索要求



## (b) 検索条件

パラグラフID

キーワード

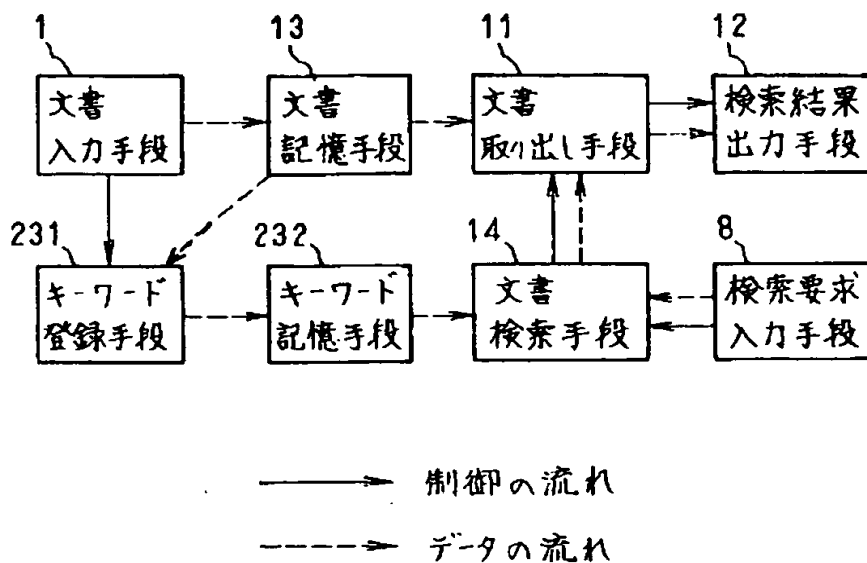
固定長/可変長フラグ

② = 『文書検索装置』 ①

## (c) 検索結果

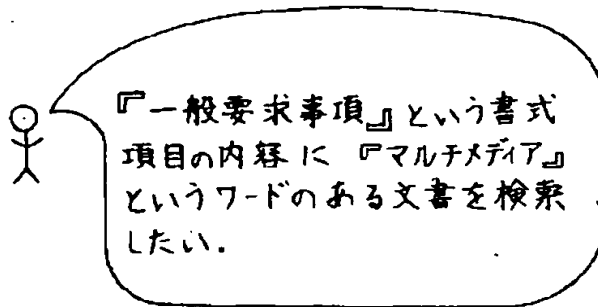
文書B

【図23】



【図18】

## (a) 検索要求



①	文書名	_____
②	適用範囲	_____
③	要求事項	_____
④	一般要求事項	<u>マルチメディア</u>
⑤	構造	_____
⑥	表示	_____

(下線部にキーワードの入力可)

## (b) 検索条件

パラグラフID

キーワード

固定長/可変長フラグ

④

=

『マルチメディア』

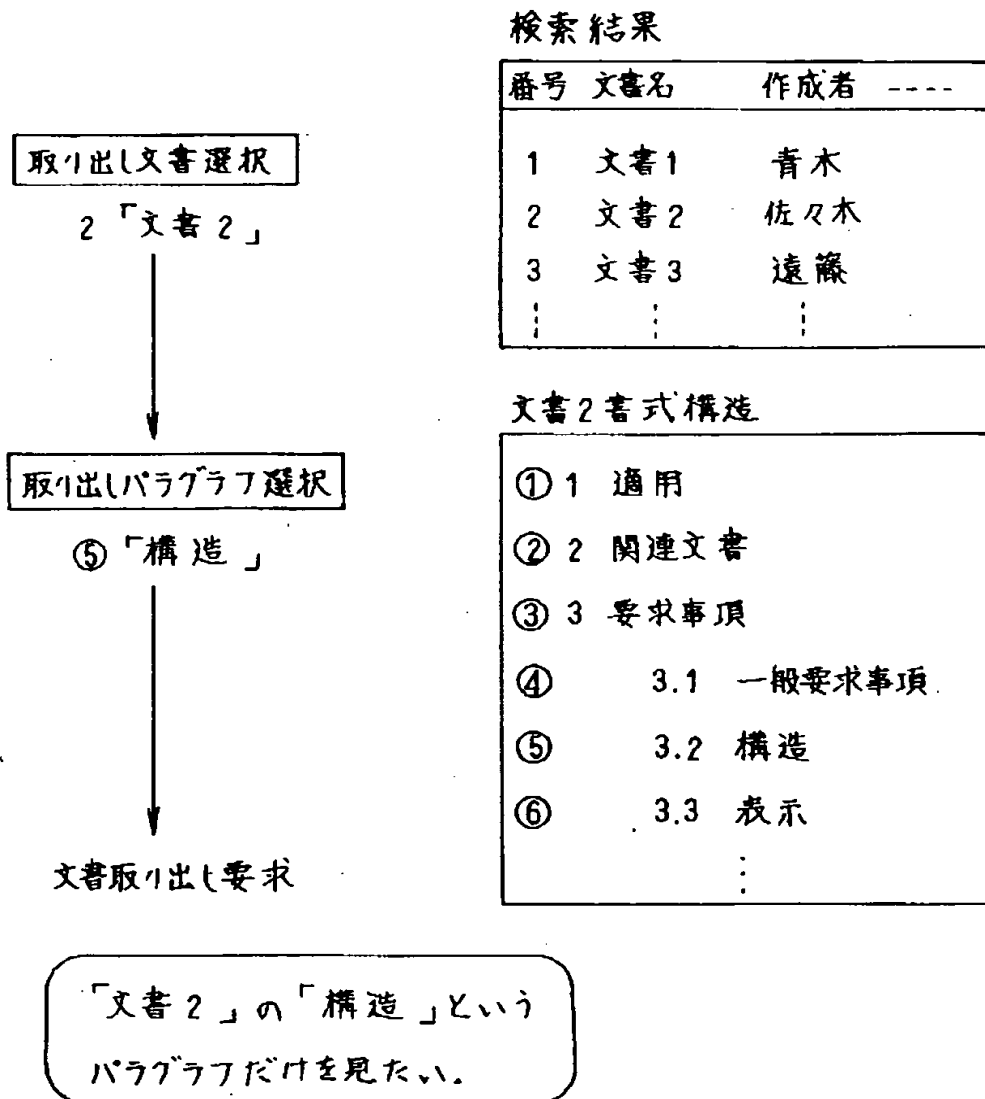
①

## (c) 検索結果

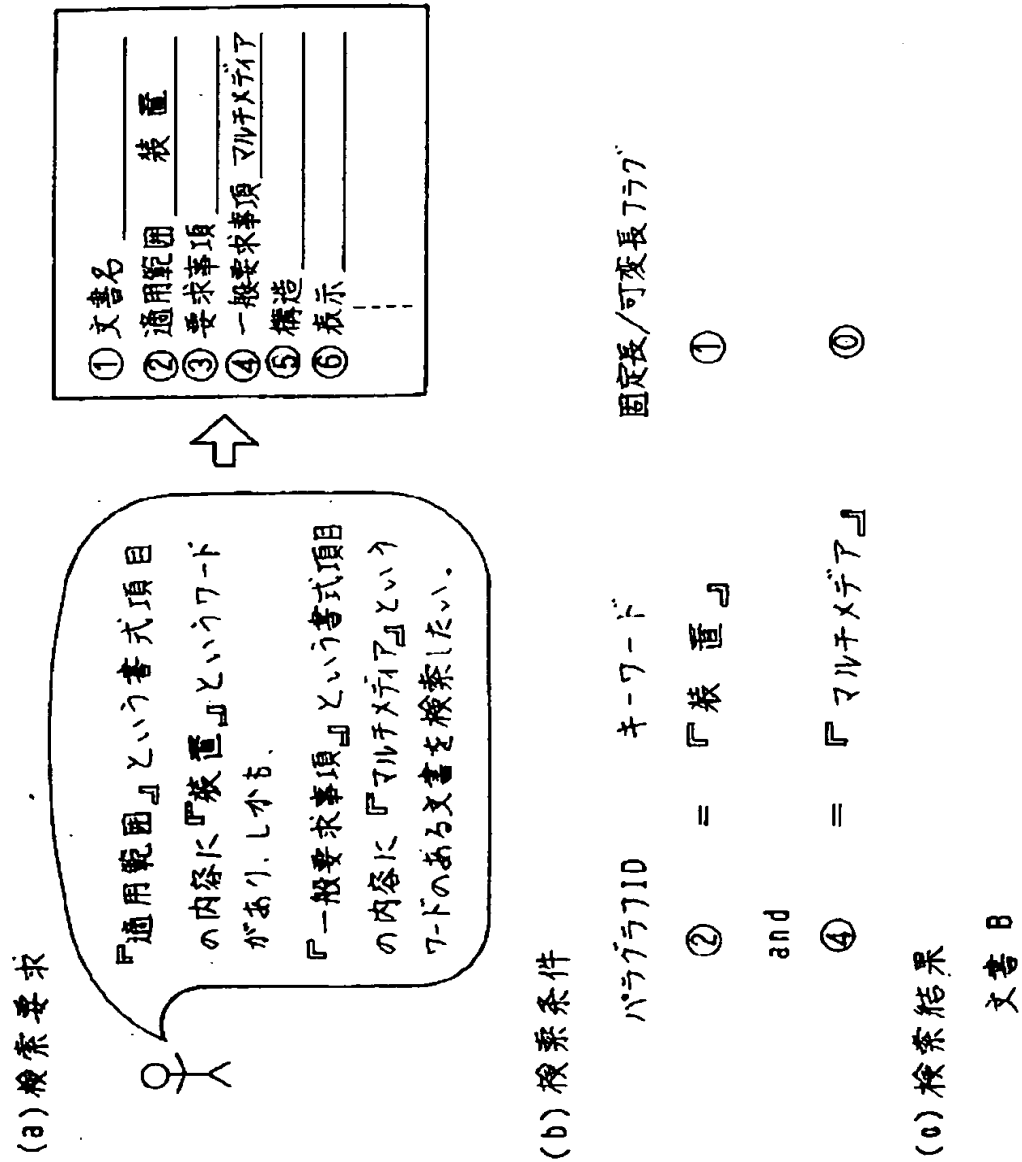
文書 A

文書 B

【図19】



【図20】





【図22】

